

Claims

What is claimed is:

- 5 1. A method in a computer system for servicing requests from one or more client computers, the method comprising:
- receiving a request from a client computer;
- a first thread processing the request by invoking a receive handler that creates a work item, wherein the first thread is part of a pool of generic threads;
- 10 a second thread performing a task specified in the work item by invoking a work handler, wherein the second thread is part of the pool of generic threads;
- receiving a result of performing the task; and
- a third thread returning at least a portion of the result to the client computer by invoking a reply handler, wherein the third thread is part of the pool of generic threads.
- 15 2. The method as claimed in claim 1, wherein receiving the request comprises:
- receiving the request into an input/output port; and
- placing a reference in a queue indicating that work is available for the first thread.
- 20 3. The method as claimed in claim 2, further comprising:
- invoking a receive thread manager when the work is available; and
- the receive thread manager scheduling the first thread for execution on one of multiple processors.
- 25 4. The method as claimed in claim 3, wherein the receive thread manager schedules the first thread from a queue of available threads.
5. The method as claimed in claim 1, further comprising the first thread

placing the work item on a work queue for execution by the second thread.

5 6. The method as claimed in claim 5, further comprising the first thread placing a reference in a completion port queue, indicating that work is available for the second thread.

10 7. The method as claimed in claim 1, further comprising the second thread creating and placing a second work item on a reply work queue when the results are received.

15 8. The method as claimed in claim 7, further comprising the second thread placing a reference in a completion port queue, indicating that work is available for the third thread.

20 9. The method as claimed in claim 1, further comprising:
receiving input data; and
the first thread storing the input data in a cache that is accessible to the second thread.

25 10. The method as claimed in claim 1, wherein the results include data, and further comprising storing the data in a cache that is accessible to the third thread.

30 11. The method as claimed in claim 1, further comprising:
determining a size of the result; and
when the size of the result is not larger than a cutoff size, storing the result in the partial results cache.

35 12. The method as claimed in claim 11, further comprising:
the third thread returning a portion of the result to client computer; and

if a second request is received for an additional portion of the result that is stored in the partial results cache, a fourth thread creating a second work item by invoking another receive handler, wherein the second work item is processed by a fifth thread, which invokes another reply handler.

5

13. The method as claimed in claim 1, further comprising:
the first thread, the second thread or the third thread indicating that the first thread, the second thread or the third thread has completed a work item; and
if a quantum has not expired for the first thread, the second thread or the third thread, then the first thread, the second thread or the third thread being given an additional work item to perform without relinquishing the central processing unit upon which the first thread, the second thread or the third thread was running.

14. A method in a computer system for servicing requests from multiple client computers, the method comprising:
receiving a request from a client computer to perform a multi-state function;
performing a first task, by a first work handler invoked by a first thread in a ready state, wherein the first task is associated with a first state of the multi-state function, and performing the first task includes issuing an asynchronous request for data;
placing the first thread back in the ready state;
receiving the data specified in the asynchronous request; and
performing a second task, by a second work handler invoked by a second thread in the ready state, wherein the second task is associated with a second state of the multi-state function, and the second task performs an operation on the data.

25

15. The method as claimed in claim 14, further comprising:
processing the request by a receive handler invoked by a third thread;
creating a work item that specifies the first task; and
placing the work item in a work queue that is accessible to the first thread.

16. The method as claimed in claim 15, wherein the first thread, the second thread, and the third thread are all identical generic threads within a pool of generic threads.

5

17. The method as claimed in claim 15, wherein placing the work item in the work queue comprises placing the work item in a low priority work queue, the method further comprising:

10 placing a second work item that specifies the second task on a high priority work queue;
a work handler looking for a work item first on the high priority work queue; and
if no work item exists on the high priority work queue, the work handler looking for the work item on the low priority queue.

15 18. The method as claimed in claim 14, wherein the asynchronous request is a request that would otherwise cause the first thread to block.

19. The method as claimed in claim 18, wherein the asynchronous request is a request that would otherwise cause a blocking condition to occur.

20

20. The method as claimed in claim 14, wherein issuing the asynchronous request comprises issuing the asynchronous request to a database manager.

21. The method as claimed in claim 20, further comprising:
25 the database manager placing the asynchronous request on a pending queue; and
when the data is received, placing a work item associated with the second state on a work queue.

22. The method as claimed in claim 21, wherein the work queue is a high

priority work queue, the method further comprising:

placing a first work item that specifies the first task on a low priority work queue;
a work handler looking for a work item first on the high priority work queue; and
if no work item exists on the high priority work queue, the work handler looking
5 for the work item on the low priority work queue.

23. The method as claimed in claim 14, further comprising returning a result
of the first task and the second task by invoking third thread, which in turn invokes a
reply handler to return the result to the client computer.

10 24. The method as claimed in claim 23, wherein the first thread, the second
thread, and the third thread are all identical generic threads within a pool of generic
threads.

15 25. The method as claimed in claim 14, further comprising:
performing additional tasks associated with subsequent function states by
additional work handlers invoked by one or more additional threads;
at least some of the additional work handlers issuing additional asynchronous
requests; and

20 placing threads associated with the at least some of the additional work handlers
back in the ready state after issuing the additional asynchronous requests.

26. A method in a computer system for servicing requests from multiple client
computers, the method comprising:

25 determining that work is available after receiving a request from a client
computer;

when work is available, a first work handler invoked by a first thread looking in a
first work queue for a first work item corresponding to the work; and

if the first work item is not found in the first work queue, the first work handler

looking in a second work queue for the first work item.

27. The method as claimed in claim 26, further comprising:
receiving a request from a client computer to perform a task;
5 creating the first work item that specifies the task;
placing the first work item in the second queue; and
indicating that the work is available.

28. The method as claimed in claim 26, further comprising:
10 the first work handler performing a task specified in the first work item; and
issuing an asynchronous request for data.

29. The method as claimed in claim 28, further comprising:
receiving the data;
15 placing a second work item on the first work queue; and
indicating that additional work is available.

30. The method as claimed in claim 29, further comprising:
a second work handler invoked by a second thread looking in the first work queue
20 for the second work item when the second work is available; and
performing a second task specified in the second work item.

31. The method as claimed in claim 26, wherein the computer system includes
multiple work queues, including the first work queue and the second work queue, each of
25 the multiple work queues are associated with a priority level, and wherein the method
further comprises:

looking for the first work item first in a work queue associated with a highest
priority level; and

if the first work item is not found in the work queue associated with the highest

priority level, looking for the first work item in each of the multiple work queues in descending priority order until the first work item is found.

32. The method as claimed in claim 31, wherein the request from the client
5 computer is a request to perform a function having multiple states, wherein each state of the multi-state function is performed by a work handler invoked by a subsequent thread based on work items placed in the multiple work queues, and wherein the work items are placed in higher and higher priority work queues as execution of the multi-state function progresses through the multiple states.

10 33. A method in a computer system for servicing requests from multiple client computers, the method comprising:

receiving, from a client computer, a request to perform a first task;
evaluating the first task, by a first handler invoked by a first thread, to determine
15 whether the first task includes complex or long-running logic; and
if the first task includes complex or long-running logic, performing the first task by a second handler invoked by a second thread.

20 34. The method as claimed in claim 33, wherein the first thread is a thread within a first group of threads having a first priority level, and the second thread is a thread within a second group of threads having a second priority level that is lower than the first priority level.

25 35. The method as claimed in claim 34, further comprising:
receiving a request to perform a second task;
evaluating the second task to determine whether the second task includes complex or long-running logic; and
if the second task does not include complex or long-running logic and a processor is not available for performing the second task, preempting the second thread and

performing the second task by a third thread in the first group of threads.

36. The method as claimed in claim 33, further comprising:
the second handler returning a result of the first task; and
5 using the result, performing a second task by a third handler invoked by a third
thread of the first group of threads.

37. The method as claimed in claim 36, wherein the first task is specified in a
first work item and the second task is specified in a second work item, the method further
10 comprising:
the first handler obtaining the first work item from a first work queue; and
the third handler obtaining the second work item from a second work queue.

38. A method in a computer system for servicing requests from one or more
15 client computers, the method comprising:
maintaining a pool of threads, wherein each thread in the pool of threads is
identical and can invoke at least one receive handler, at least one work handler, and at
least one reply handler;
invoking receive handlers by the threads in response to receiving requests from
20 one or more client computers, wherein the receive handlers create work items that specify
tasks to be performed to satisfy the request;
invoking work handlers by the threads to perform the tasks specified in the work
items;
receiving results of the tasks; and
25 invoking reply handlers by the threads to return at least portions of the results to
the client computers.

39. The method as claimed in claim 38, wherein one or more of the receive
handlers, one or more of the work handlers, and one or more of the reply handlers can be

simultaneously executed on multiple processors available to the computer system.

40. The method as claimed in claim 38, wherein the pool of threads includes a number of threads in a range from about $N+1$ to about $2*N$, where N is a number of
5 processors available to the computer system.

41. A method in a computer system for servicing requests from multiple client computers, the method comprising:
monitoring a quantity of work being performed by the computer system;
10 determining whether the quantity has exceeded an upper limit; and
if the quantity has exceeded the upper limit but has not dropped below a lower limit, not accepting new requests into the computer system.

42. The method as claimed in claim 41, further comprising, if the quantity has
15 exceeded the upper limit and has dropped below the lower limit, accepting the new requests into the computer system.

43. The method as claimed in claim 41, further comprising continuing to monitor the quantity of work being performed by the computer system after the quantity
20 has exceeded the upper limit.

44. The method as claimed in claim 41, wherein the quantity of work is indicated by a number of work items on one or more work queues, and determining whether the quantity has exceeded the upper limit comprises determining whether the
25 number of work items has exceeded a predefined value.

45. A method in a computer system for servicing requests from multiple client computers, the method comprising:
monitoring an amount of time to return a result by the computer system;

determining whether the amount of time has exceeded an upper limit; and
if the amount of time has exceeded the upper limit but has not dropped below a lower limit, not processing new work items.

5 46. The method as claimed in claim 45, further comprising, if the amount of time has exceeded the upper limit and has dropped below the lower limit, accepting and processing the new work items.

10 47. The method as claimed in claim 45, further comprising continuing to monitor the amount of time after the amount of time has exceeded the upper limit.

15 48. The method as claimed in claim 45, wherein the amount of time includes a time it takes to post the result to an output port, and determining whether the amount of time has exceeded the upper limit comprises determining whether the time it takes to post the result has exceeded a predefined value.

20 49. An application program for implementation by an application server, the application program comprising:
 at least one receiver handler that can be invoked by a thread within a pool of threads;
 at least one work handler that also can be invoked by the thread; and
 at least one reply handler that also can be invoked by the thread.

25 50. The application program as claimed in claim 49, wherein a receiver handler of the at least one receiver handler is executed when the application server receives a request from a client computer, and the receiver handler creates a first work item to be performed by a work handler of the at least one work handler.

51. The application program as claimed in claim 50, wherein the work handler

is executed when the first work item exists, the work handler receives results, and the work handler creates a second work item to be performed by a reply handler of the at least one reply handler.

5 52. The application program as claimed in claim 51, wherein the reply handler is executed when the second work item exists, and the reply handler sends the results to the client computer.

10 53. The application program as claimed in claim 49, wherein the thread can invoke multiple work handlers, where some of the multiple work handlers are designed to perform tasks associated with various states of a multi-state function, and at least some of the multiple work handlers issue asynchronous requests for data when a state transition is to be performed, wherein the thread is then placed back in a ready state to execute a subsequent work item.

15 54. The application program as claimed in claim 53, wherein when the data is returned, a second work handler is executed.

20 55. The application program as claimed in claim 49, wherein multiple identical copies of the thread exist within the pool of threads, and at least one copy of the thread is executed each time a request is received from a client computer.

25 56. The application program as claimed in claim 49, wherein multiple copies of the thread can simultaneously be executed by multiple processors available to the application server.

 57. The application program as claimed in claim 49, further comprising one or more complex logic handlers that can be invoked by a second type of thread, wherein a thread of the second type of thread is executed when a request from a client computer

involves execution of complex or long-running logic.

58. A computer system for servicing requests from one or more client computers, the computer system comprising:

- 5 a memory containing an application server which maintains a pool of threads, wherein each thread in the pool of threads is identical and can invoke at least one receive handler, at least one work handler, and at least one reply handler, and the application server further schedules threads in the pool of threads for execution on multiple processors available to the computer system; and
- 10 the multiple processors for simultaneously executing multiple threads within the pool of threads.

59. The computer system as claimed in claim 58, wherein the application server further performs functions of:

- 15 receiving a request from a client computer;
a first thread processing the request by invoking a receive handler which creates a work item, wherein the first thread is part of the pool of threads;
a second thread performing a task specified in the work item by invoking a work handler, wherein the second thread is part of the pool of threads;
- 20 receiving a result of performing the task; and
a third thread returning at least a portion of the result to the client computer by invoking a reply handler, wherein the third thread is part of the pool of threads.

60. The computer system as claimed in claim 58, wherein the application server further performs functions of:

- 25 receiving a request from a client computer to perform a multi-state function;
a first thread within the pool of threads performing a first task by invoking a first work handler, wherein the first task is associated with a first state of the multi-state function, and performing the first task includes issuing an asynchronous request for data;

placing the first thread back in a ready state;
receiving the data specified in the asynchronous request; and
a second thread within the pool of threads performing a second task by invoking a
second work handler, wherein the second task is associated with a second state of the
5 multi-state function, and the second task performs an operation on the data.

61. The computer system as claimed in claim 58, wherein the application
server further performs functions of:
determining that work is available after receiving a request from a client
10 computer;
when work is available, a first thread within the pool of threads invoking a first
work handler to look in a first work queue for a first work item corresponding to the
work; and
if the first work item is not found in the first work queue, the first work handler
15 looking in a second work queue for the first work item.

62. The computer system as claimed in claim 58, wherein the application
server further performs functions of:
receiving, from a client computer, a request to perform a first task;
20 evaluating the first task, by a first handler invoked by a first thread within the pool
of threads, to determine whether the first task includes complex or long-running logic;
and

if the first task includes complex or long-running logic, performing the first task
by a second handler invoked by a second thread that is not within the pool of threads.
25

63. The computer system as claimed in claim 58, wherein the application
server further performs functions of:
monitoring a quantity of work being performed by the computer system;
determining whether the quantity has exceeded an upper limit; and

if the quantity has exceeded the upper limit but has not dropped below a lower limit, not accepting new requests into the computer system.

5 64. The computer system as claimed in claim 58, wherein the application server further performs functions of:
 monitoring an amount of time to return a result by the computer system;
 determining whether the amount of time has exceeded an upper limit; and
 if the amount of time has exceeded the upper limit but has not dropped below a lower limit, not processing new work items.

10 65. A computer-readable medium holding computer executable instructions, the computer-readable medium for performing a method in a computer system for, the method comprising:
 maintaining a pool of threads, wherein each thread in the pool of threads is
15 identical and can invoke at least one receive handler, at least one work handler, and at least one reply handler, and the application server further schedules threads in the pool of threads for execution on multiple processors available to the computer system; and
 simultaneously executing multiple threads within the pool of threads on multiple processors available to the computer system.

20 66. The computer-readable medium as claimed in claim 65, wherein the method further comprises:
 receiving a request from a client computer;
 a first thread processing the request by invoking a receive handler, which creates a
25 work item, wherein the first thread is part of the pool of threads;
 a second thread performing a task specified in the work item by invoking a work handler, wherein the second thread is part of the pool of threads;
 receiving a result of performing the task; and
 a third thread returning at least a portion of the result to the client computer by

invoking a reply handler, wherein the third thread is part of the pool of threads.

67. The computer-readable medium as claimed in claim 65, wherein the method further comprises:

- 5 receiving a request from a client computer to perform a multi-state function;
a first thread within the pool of threads performing a first task by invoking a first work handler, wherein the first task is associated with a first state of the multi-state function, and performing the first task includes issuing an asynchronous request for data;
placing the first thread back in a ready state;
- 10 receiving the data specified in the asynchronous request; and
a second thread within the pool of threads performing a second task by invoking a second work handler, wherein the second task is associated with a second state of the multi-state function, and the second task performs an operation on the data.

15 68. The computer-readable medium as claimed in claim 65, wherein the method further comprises:

- determining that work is available after receiving a request from a client computer;
- when work is available, a first work handler invoked by a first thread within the
- 20 pool of threads looking in a first work queue for a first work item corresponding to the work; and
- if the first work item is not found in the first work queue, the first work handler looking in a second work queue for the first work item.

25 69. The computer-readable medium as claimed in claim 65, wherein the method further comprises:

- receiving, from a client computer, a request to perform a first task;
- evaluating the first task, by a first handler invoked by a first thread within the pool of threads, to determine whether the first task includes complex or long-running logic;

and

if the first task includes complex or long-running logic, performing the first task by a second handler invoked by a second thread that is not within the pool of threads.

- 5 70. The computer-readable medium as claimed in claim 65, wherein the method further comprises:
- monitoring a quantity of work being performed by the computer system;
 - determining whether the quantity has exceeded an upper limit; and
 - if the quantity has exceeded the upper limit but has not dropped below a lower
- 10 limit, not accepting new requests into the computer system.

71. The computer-readable medium as claimed in claim 65, wherein the method further comprises:
- monitoring an amount of time to return a result by the computer system;
 - 15 determining whether the amount of time has exceeded an upper limit; and
 - if the amount of time has exceeded the upper limit but has not dropped below a lower limit, not processing new work items.